



CentraleSupélec

# AI Challenge

Presented by :

**Tudo Bein**

**Members of the group :**

- João Regazzi
- Guilherme Mertens
- Lucas Tramonte
- Arthur Vogels
- Rebecca Bayssari

**Presented to :**

- Geraud Faye
- Jean-Philippe Poli

# Context

1- Dataset Analysis

2- Overview

3- Pipeline

4- Models

5- Evaluation and results

6- Improvements

# Dataset Analysis

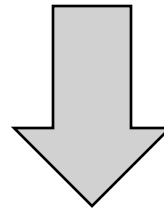
## DAM training Dataset

- Limited to one image per class
- Uniform Background: Mostly white backgrounds



## Test image Dataset

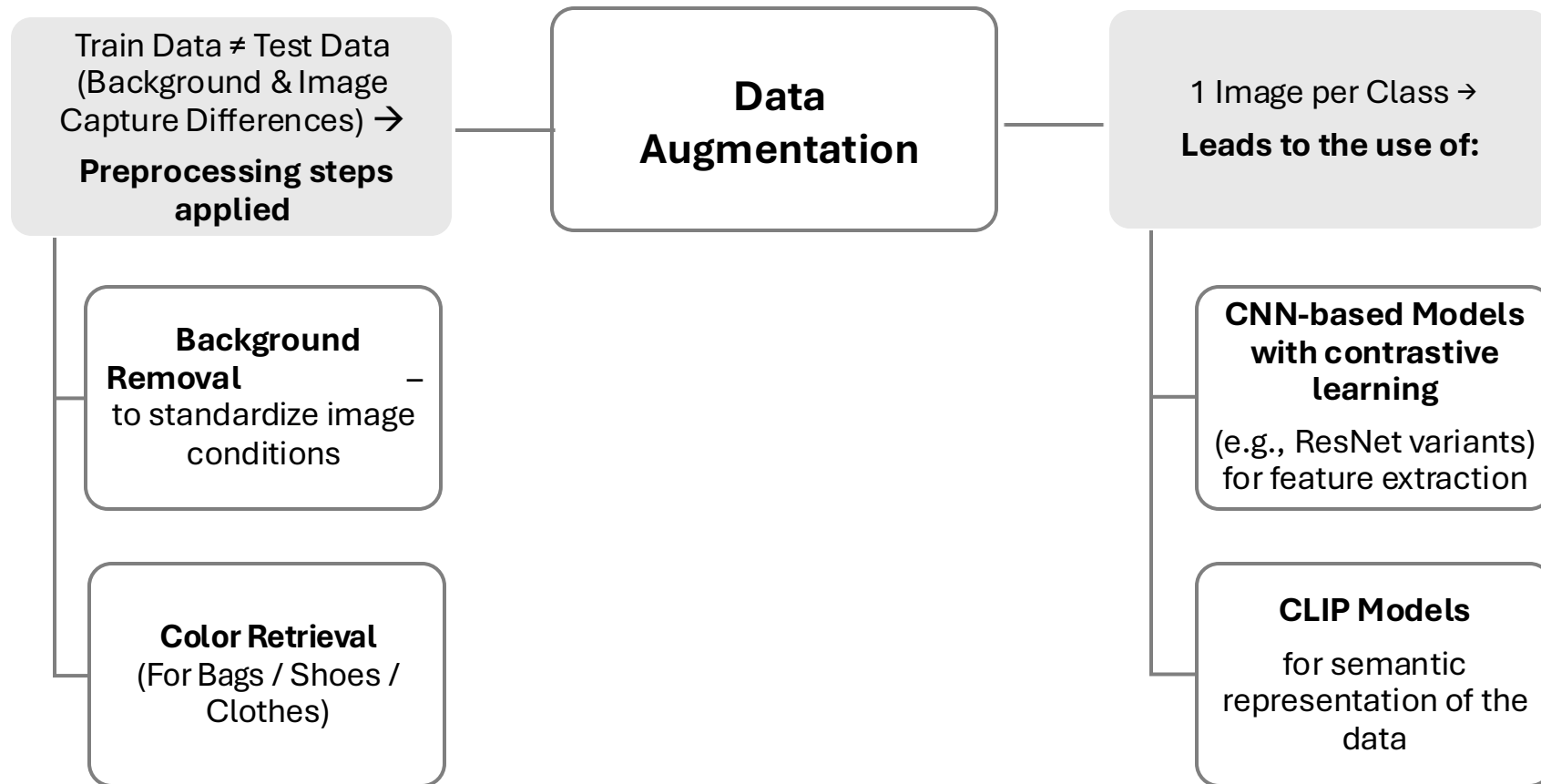
- Diverse backgrounds (different places, stores)
- Real-world variability and complexity



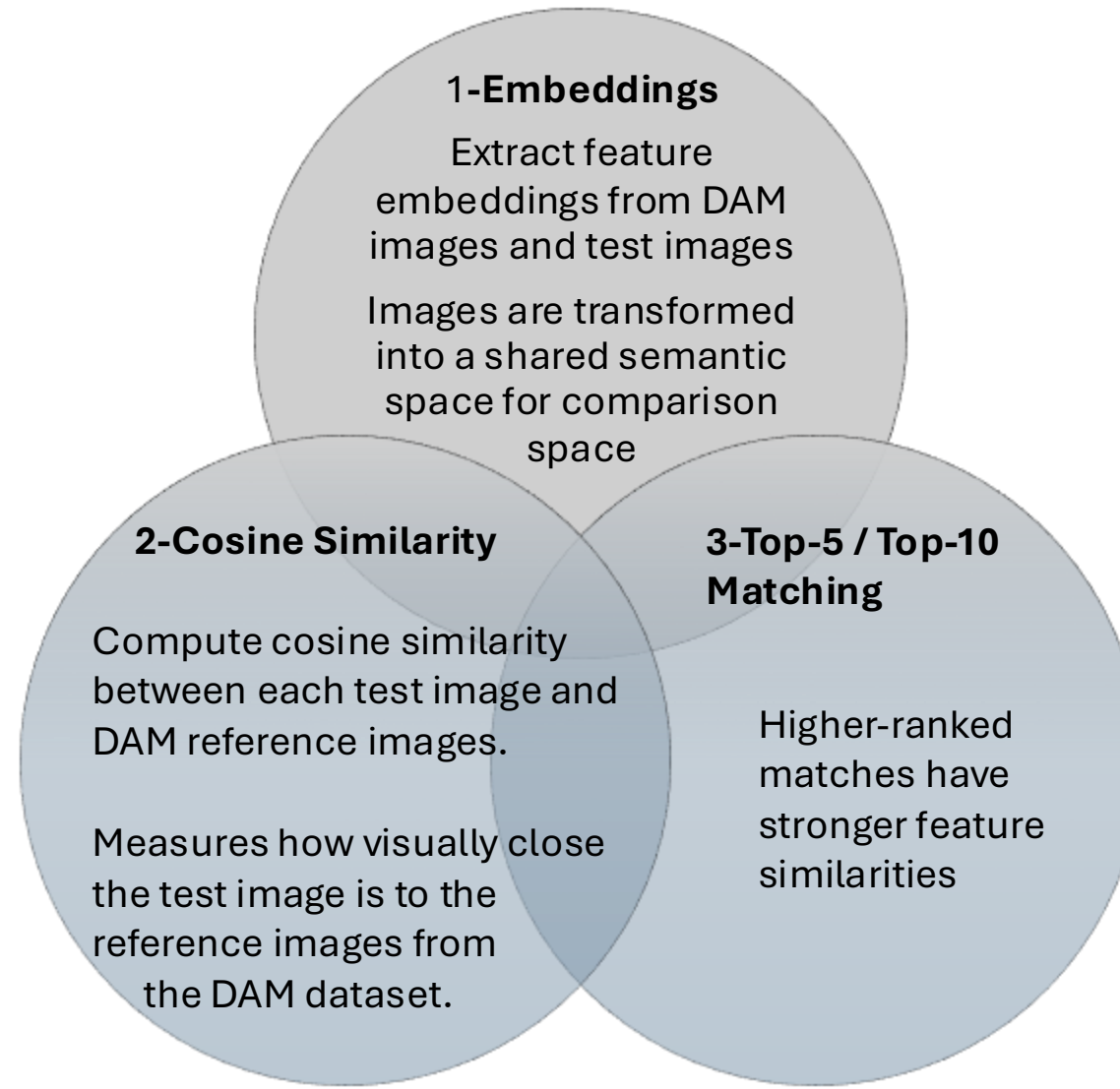
## One shot Learning Challenge

- Definition: Learning from just one example per class
- Approach: Train a model to learn semantic representations of the images and compare them with similarity measures

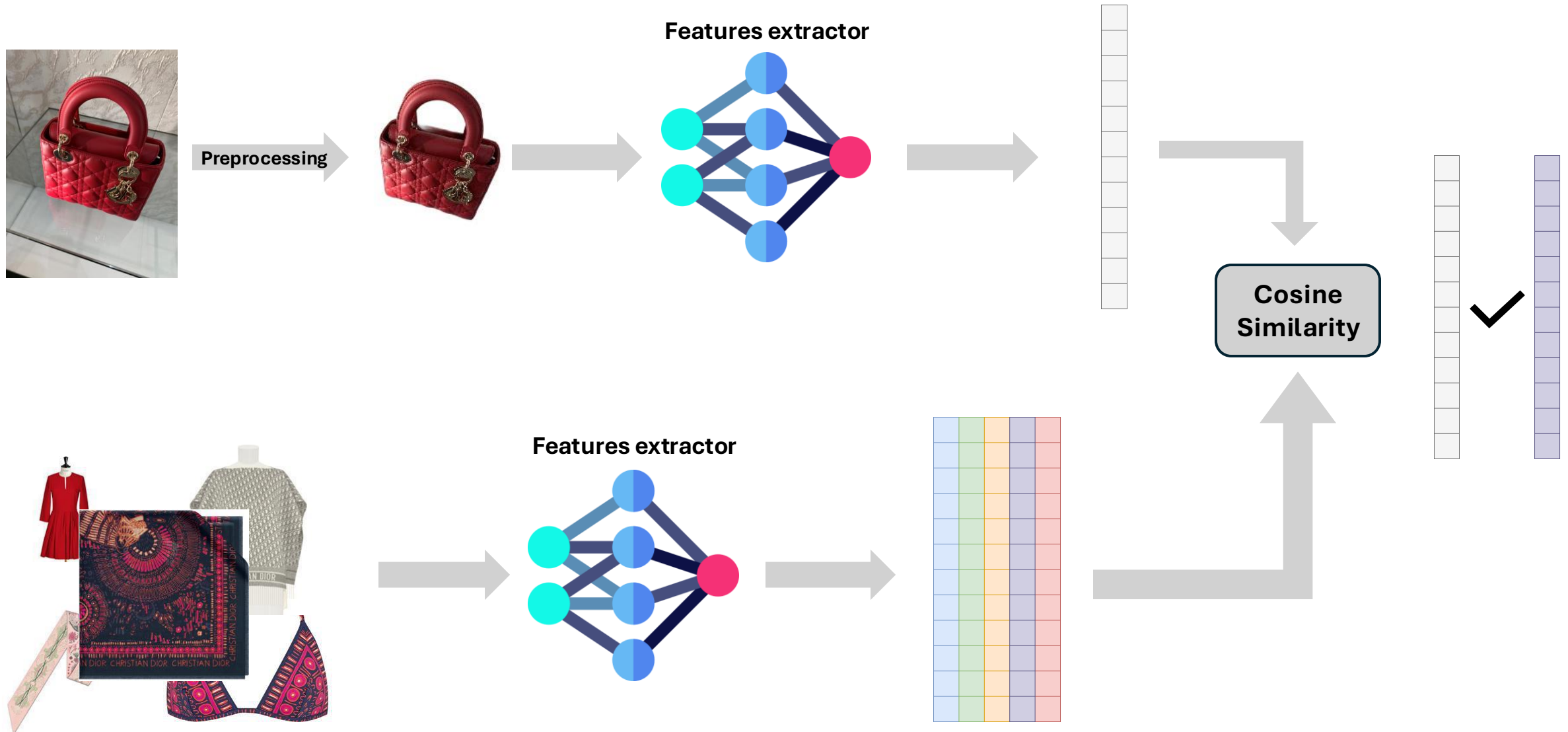
# Problem Overview



# Pipeline



# Pipeline



## 1. Image Embedder using pretrained Resnet Backbones

- Test the pipeline with a simple, pretrained model
- Use of Resnet50 pretrained on ImageNet with no finetuning

Those embeddings are not optimized for classification of unseen classes, leading to not so great results

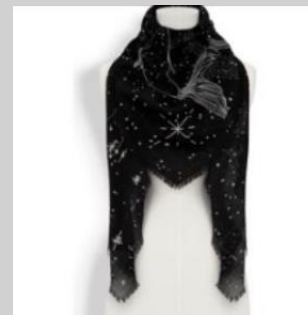


### 2. Finetuning ResNet with contrastive learning

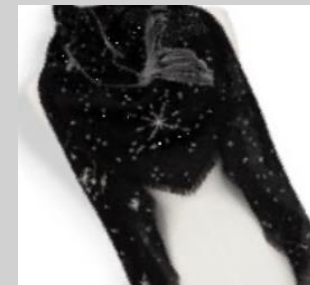
**Contrastive learning:** Different images of a same object should have similar embeddings, while images of different objects should have dissimilar embeddings.



**We compare original images with augmented ones:**



=

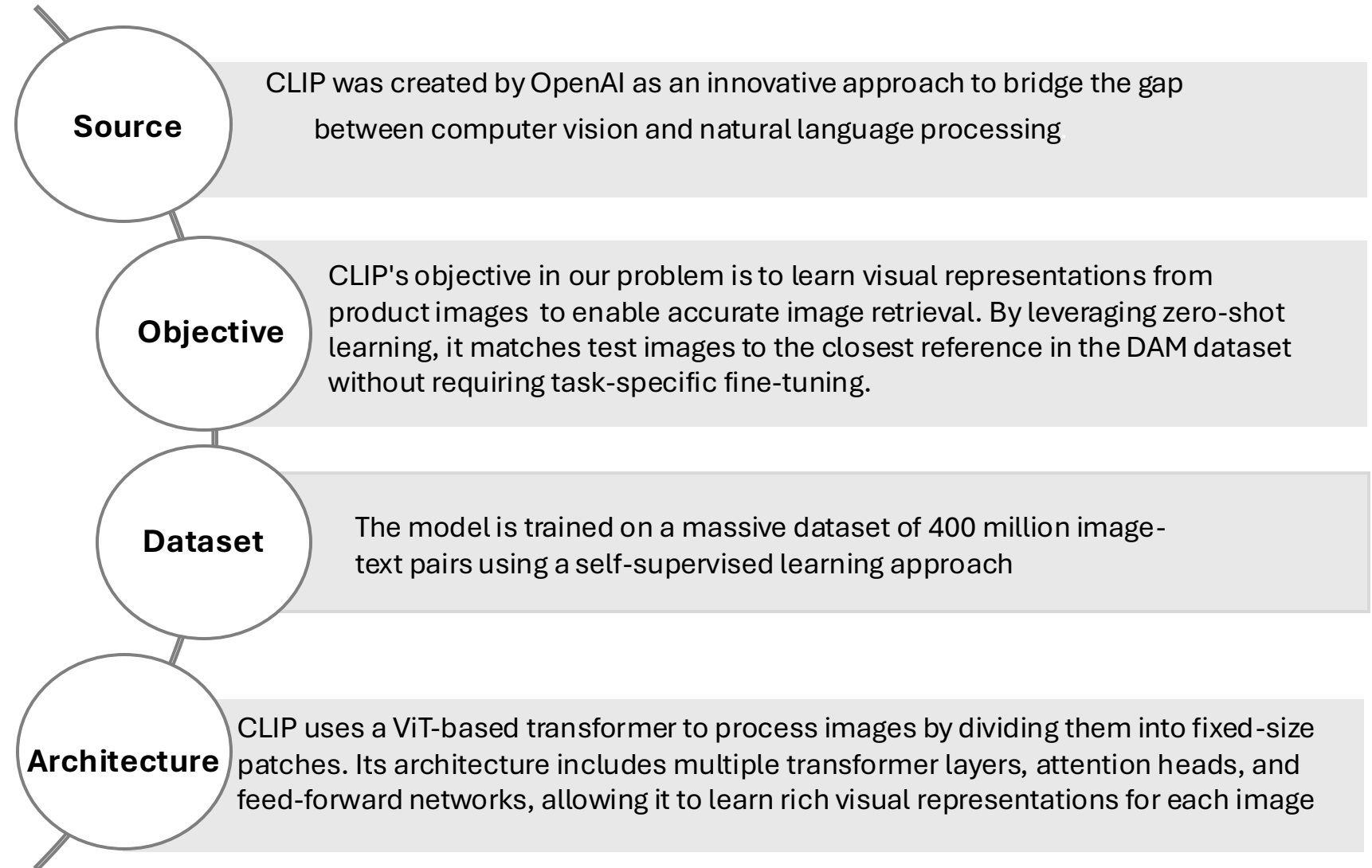


≠





# CLIP



# The Final approach

In particular, we use a fine-tuned version of CLIP called:

## **Fashion Clip\***

### **WHY**

- FashionCLIP is designed for real-world applications and is openly available
- It is fine-tuned for the fashion industry, making it highly relevant to our challenge

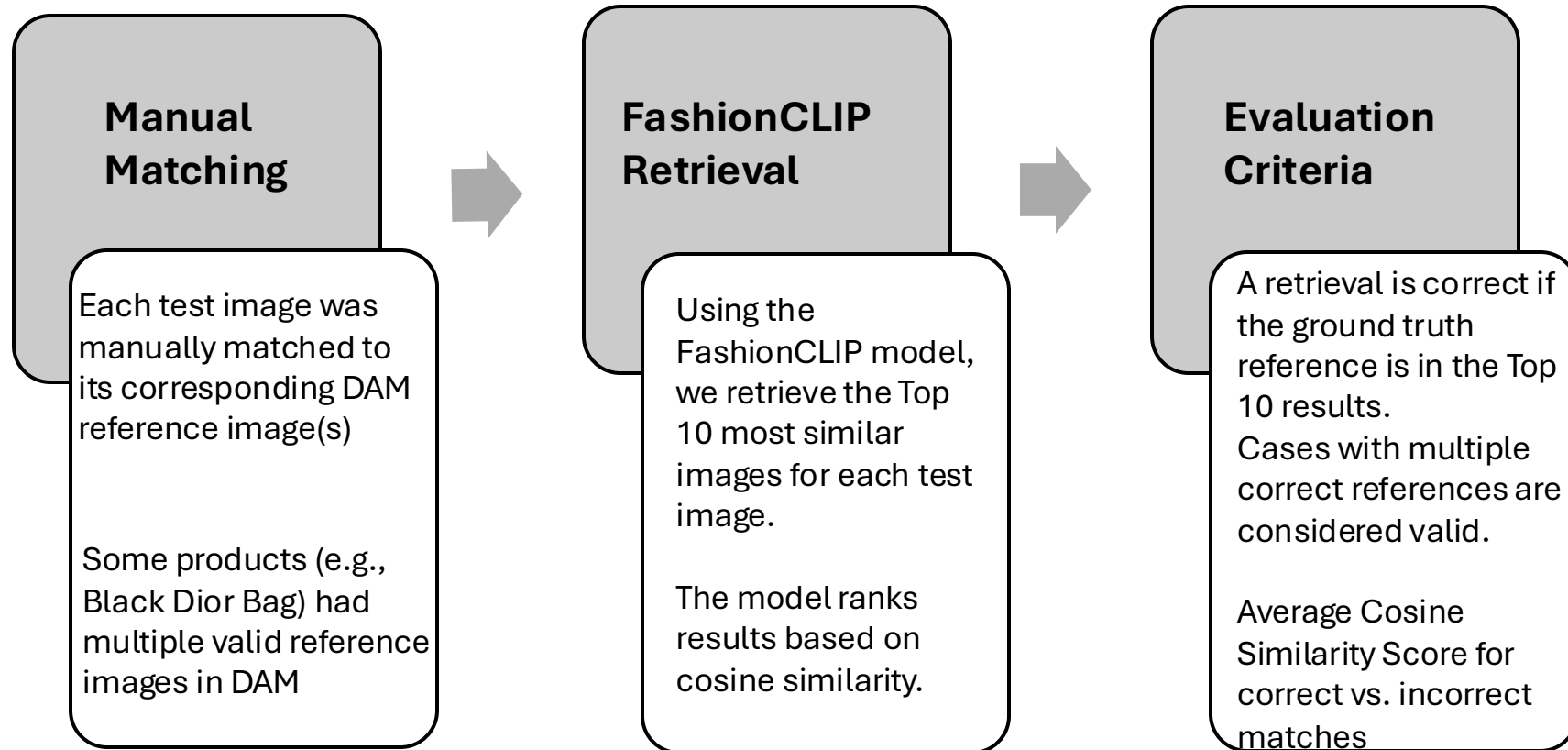
### **Architecture**

- Uses ViT-B/32 (Vision Transformer) for image encoding.

### **Training**

- Trained on 700k+ fashion image-text pairs (luxury retailers, online stores).
- Fine-tuned to recognize nuanced fashion concepts with zero-shot learning.

\* <https://github.com/patrickjohncyh/fashion-clip>



# Results

Model	Top 1 accu.	Top 5 accu.	Top 10 accu.
ResNet50	21%	45%	57%
Clip	21%	47%	58%
Fashion-Clip	<b>47%</b>	<b>83%</b>	<b>91%</b>

Comparison of the different models

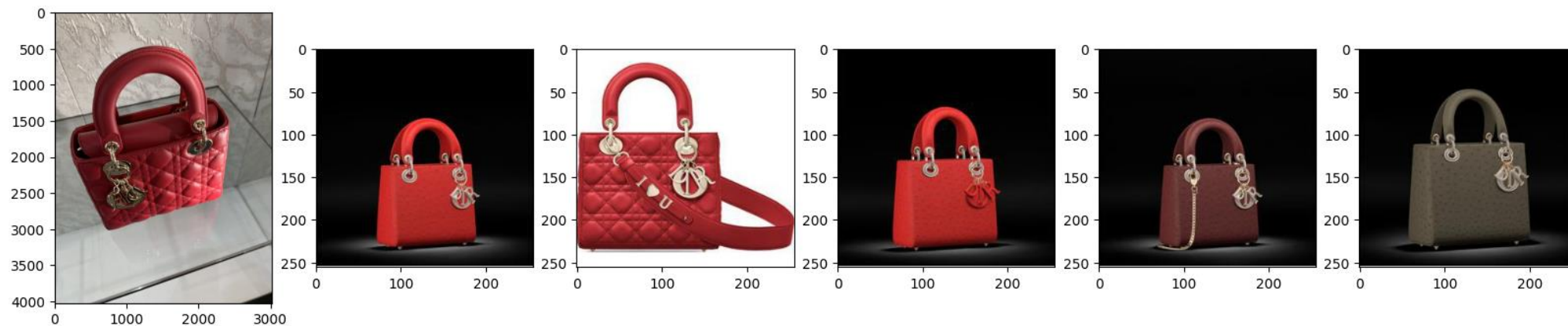


Qualitative results of the FashionClip model

# Results



Qualitative results of the ResNet model

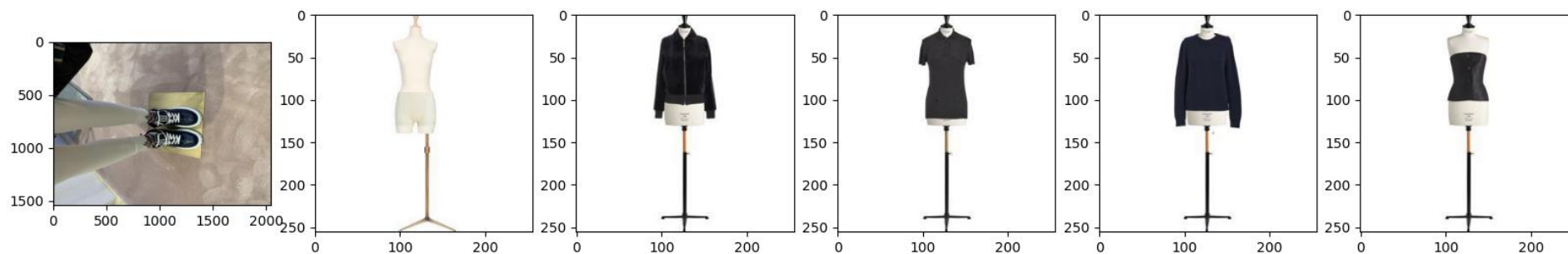


Qualitative results of the Clip model

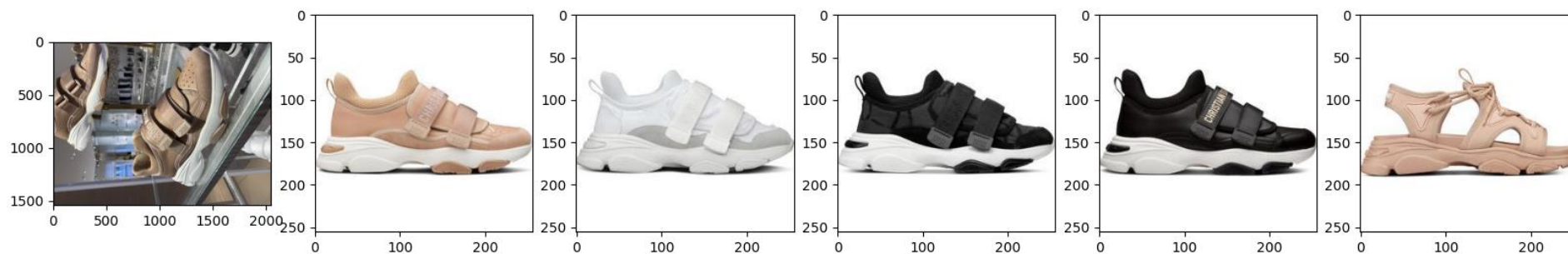
# Results

Fashion Clip works well but is not perfect!

*Class hallucinations:*



*Color hallucinations:*



# Improvements

Fashion Clip sometimes misunderstand colors, leading to the idea of adding a color analysis to our pipeline:

If  $Distance(A_{colour}, B_{colour}) < threshold$ :

$$Similarity(A, B) = \cos(A_{emb}, B_{emb}) + \alpha (1 - Distance(A_{colour}, B_{colour}) / threshold)$$

Else:

$$Similarity(A, B) = \cos(A_{emb}, B_{emb})$$



Model	Top 1 accu.	Top 5 accu.	Top 10 accu.
ResNet50	21%	45%	57%
Clip	21%	47%	58%
Fashion-Clip	47%	83%	91%
Fashion-Clip + Color	<b>51%</b>	<b>84%</b>	<b>94%</b>

The best model was :

## **Fashion Clip with enhanced color analysis**

### **Achievements**

- Successfully implemented CLIP for effective image retrieval, leveraging zero-shot learning for unseen classes.
- Achieved over 94% accuracy in top 10 accuracy for the best model.

### **Future Work**

- Experiment with ensemble methods to combine multiple models for improved retrieval accuracy.
- Explore additional augmentation strategies to simulate real-world conditions during training.



# References

[https://huggingface.co/docs/transformers/en/model\\_doc/owlv2](https://huggingface.co/docs/transformers/en/model_doc/owlv2)

[https://github.com/NielsRogge/Transformers-Tutorials/blob/master/OWLv2/Zero\\_and\\_one\\_shot\\_object\\_detection\\_with\\_OWLv2.ipynb](https://github.com/NielsRogge/Transformers-Tutorials/blob/master/OWLv2/Zero_and_one_shot_object_detection_with_OWLv2.ipynb)

<https://pypi.org/project/rembg/>

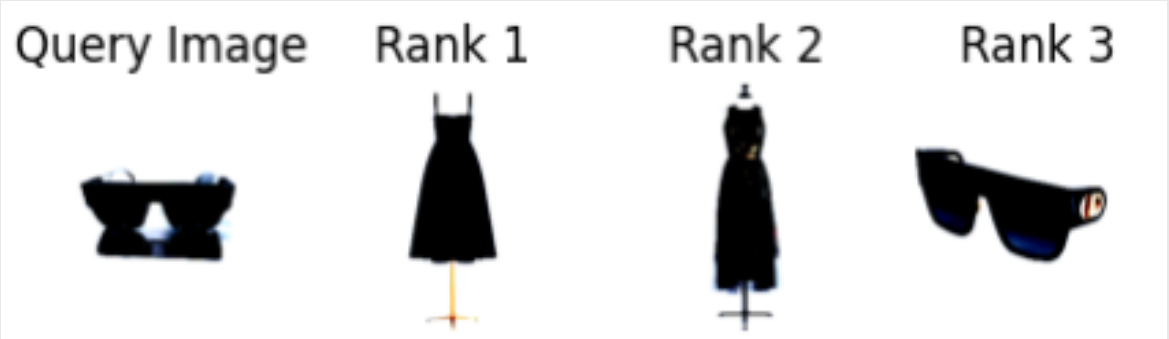
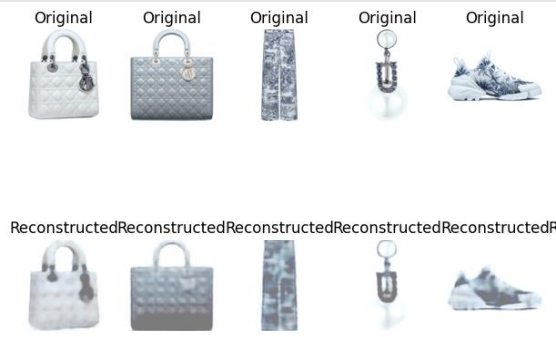
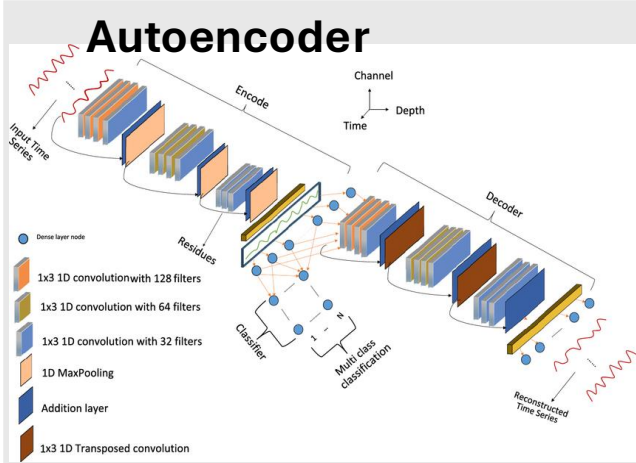
<https://huggingface.co/patrickjohnncyh/fashion-clip>

[https://youtu.be/oEKg\\_jiV1Ng?si=\\_8l9pBSq6BseF86l](https://youtu.be/oEKg_jiV1Ng?si=_8l9pBSq6BseF86l)

Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., ... & Sutskever, I. (2021, July). Learning transferable visual models from natural language supervision. In International conference on machine learning (pp. 8748-8763). PMLR.

Melekhov, I., Kannala, J., & Rahtu, E. (2016, December). Siamese network features for image matching. In 2016 23rd international conference on pattern recognition (ICPR) (pp. 378-383). IEEE.

# Exploring other Models for Image retrieval (ANNEXE)



## OwlViT (Owlv2)

It's a vision transformer model designed for open-world object detection and image retrieval, specifically the "google/owlvit-base-patch32" variant. The model extracts image features (embeddings) using a pretrained transformer architecture, and cosine similarity is then used to match test images to DAM reference images based on visual similarity.



# Data Augmentation



01CDO370I001C810



01CDO370I001C810\_aug\_0



01CDO370I001C810\_aug\_1



01CDO370I001C810\_aug\_2



01CDO370I001C810\_aug\_3



01CDO370I001C810\_aug\_4



01CDO370I001C810\_aug\_5



01CDO370I001C810\_aug\_6

